

Speech Technology For Supporting Community-Based Endangered Language Documentation

Robbie Jimerson^{1,2}, Richard Hatcher³, Raymond Ptucha², Emily Prud’hommeaux^{2,4}

¹Seneca Nation of Indians, ²Rochester Institute of Technology

³State University of New York at Buffalo, ⁴Boston College



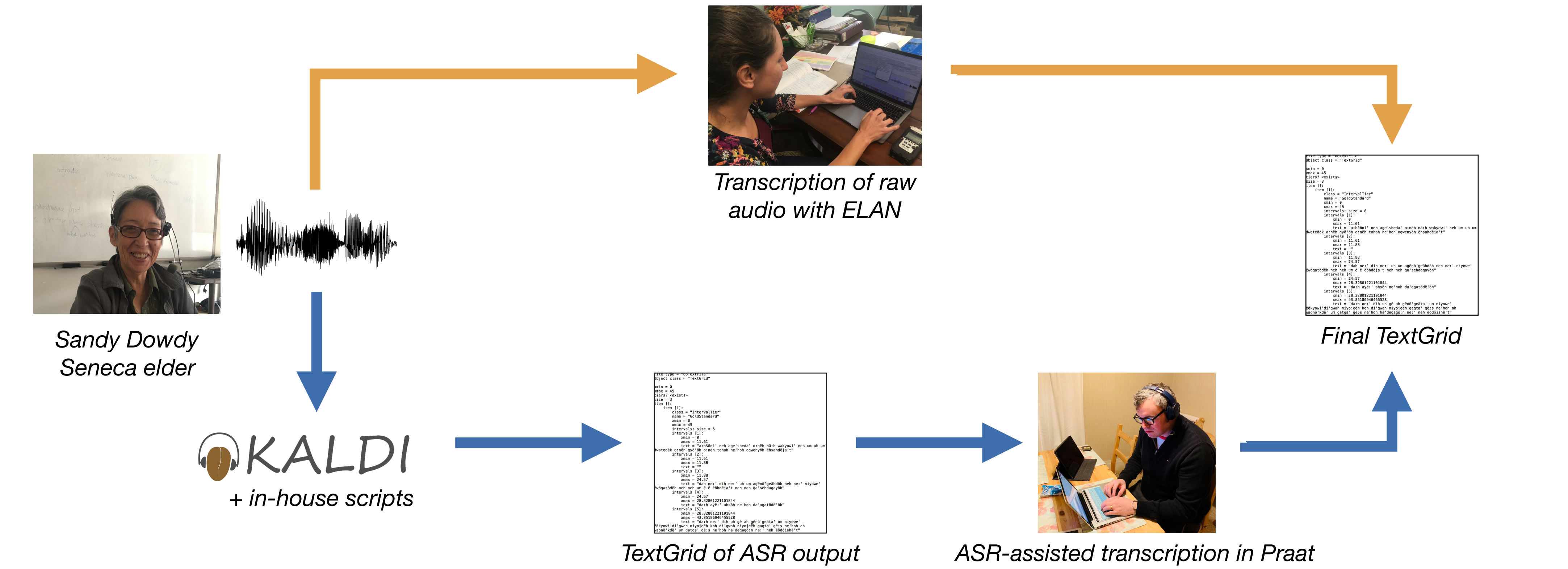
PROJECT GOALS

- Explore **automatic speech recognition (ASR)** for supporting transcription of **Seneca**
- Evaluate the **speed and accuracy** of ASR-supported transcription as a function of language skill.
- Collect data on **individual preferences** for correcting ASR output vs. producing transcripts “from scratch”.

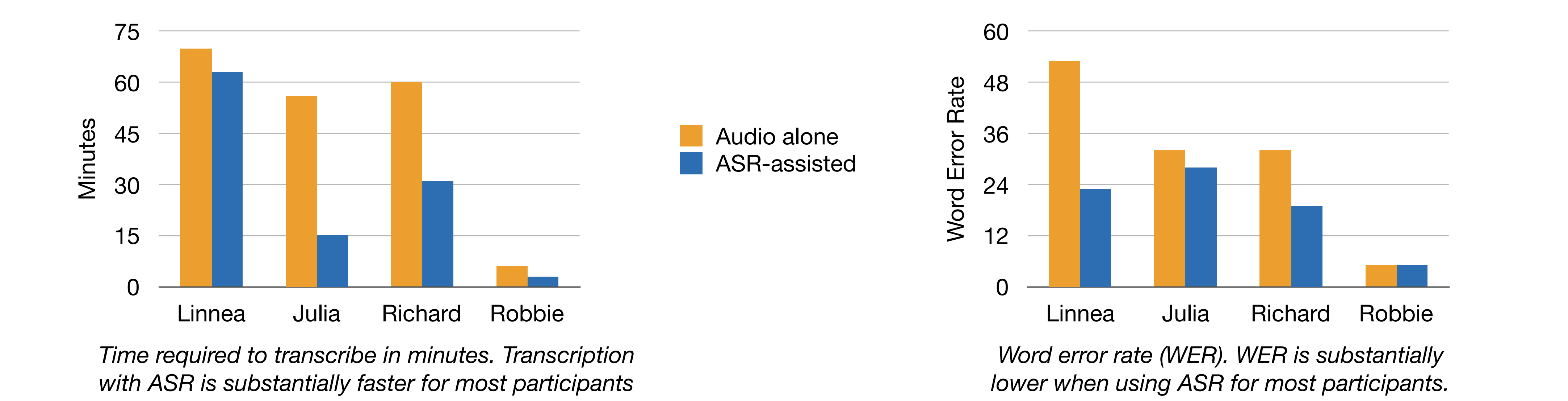
ONÖDOWA’GA:’ GAWËNÖ’

- Seneca**: member of the Iroquoian language family.
- ~**50** first-language speakers, ~**100+** learners.
- Spoken primarily in **Western NYS and Ontario**.
- Polysynthetic** morphology, highly agglutinative.
- Complex morphophonology**.

TRANSCRIPTION PIPELINES



USABILITY STUDY: ASR-ASSISTED TRANSCRIPTION



Linnea	undergrad linguistics RA	“I had a difficult time with the ASR, because I spent more time cross-checking the transcription than actually just transcribing.”
Julia	undergrad linguistics RA	“Using ASR, I was able to focus on comparing the audio to the transcription rather than trying to perceive what was being said.”
Richard	PhD student researching Cayuga	“I preferred the ASR because it took care of some of the necessary steps, e.g., segmenting speech into utterance units.”
Robbie	fluent Seneca L2 speaker	“The only thing that slowed me down using ASR was having to copy and paste around the incorrect words.”

ACKNOWLEDGEMENTS

This material is based upon work supported by the National Science Foundation under Grant No. 1761562.

We are grateful for the support and generosity of the elders of the Seneca Nation of Indians.